# Safety Last: Analysis of the Rayleigh Curve with Normalized Software Data

**Dan Strickland**
**Cost Analysis and Parametric Estimation,**
**Software SME**

DISTRIBUTION STATEMENT A. Approved for public release; distribution is unlimited.

# Overview

- **Rayleigh Curve**
- **Normalizing the SRDR Database**
- **SRDR Using Rayleigh**
- **SRDR Quality Issues**
- **Rayleigh Peak Staffing Date Benchmarks**
- **Rayleigh as Weibull**
- **Boehm, Stutzke, and Strickland Rayleigh Curves**
- **Staffing Profiles Using Rayleigh**
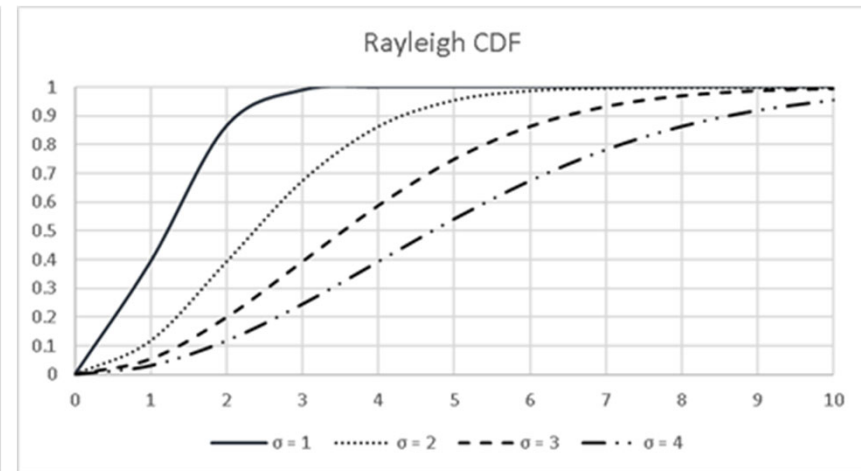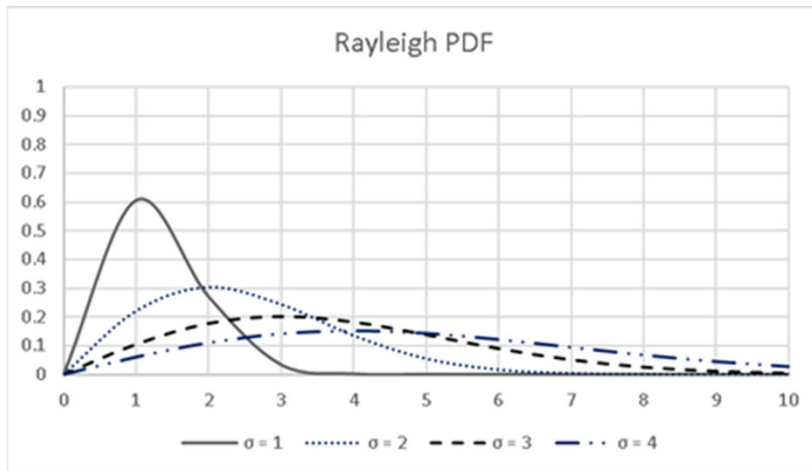- **Rayleigh and the New SRDR**
- **Future Work and Conclusions**

# *Safety Last!*

Software Estimation in DoD tends to focus more on effort and cost with little regard to time and phasing

# Rayleigh Curve

- **A continuous probability distribution named for John William Strutt, the 3rd Baron Rayleigh**
- **Defined by a positive shape parameter (σ):**
  - $f(x; \sigma) = \dfrac{x}{\sigma^2} e^{-x^2/(2\sigma^2)}, x \geq 0$ (PDF)
  - $F(x; \sigma) = 1 - e^{-x^2/(2\sigma^2)}$ (CDF)

# Rayleigh Curve - Staffing

- **Peter Norden first identified Rayleigh curves with staffing profiles – some staffing looks like Rayleigh**
- **Lawrence Putnam first applied Norden-Rayleigh curves to software staffing levels in his Software Lifecycle Model (SLIM)**
- **Norden-Rayleigh curves:**

  - $FTE(t) = \dfrac{Kt}{t_d{}^2} \, exp\left[\dfrac{-t^2}{[2t_d{}^2]}\right]$

    – FTE(t) = full-time equivalent personnel at time t
    – K = total project effort in man-months
    – $t_d$ = point in time where peak staffing occurs

  - $E(t) = K\left(1 - exp\left[\dfrac{-t^2}{(2t_d{}^2)}\right]\right)$

    – E(t) = total effort expended from 0 to time t

---

Can Rayleigh curves be used with DoD Software Resources Data Reporting (SRDR) data to develop time-phasing benchmarks?

---

# SRDR Database Ground Rules and Assumptions

- **Starting with the SRDR database of NOV 2018 (4084 records)**
  - Final SRDRs
  - "Good" Quality Tag
  - Populated Application Domain field
- **All data items should be of component or CSCI "size" in ESLOC:**
  - MDA Equivalent SLOC (ESLOC) = New + 50% (Modified) + 5% (Reuse) + 30% (AutoGen)
  - CSCI size is greater than 5K ESLOC, less than 200K ESLOC (same as Aerospace study)
- **All data items should have defined hours for Software Design, Code, and Test & Integration (DCTI)**
  - **Architecture/Design** hours are SW Design hours
  - **Code and Unit Test** hours are SW Code hours
  - **SW and System Integration, SW Qualification Testing** hours are SW Test and Integration hours
  - **Requirements Analysis** and **SW Developmental Test and Evaluation (DT&E)** hours are **not** part of DCTI hours
  - **Other** hours are distributed proportionally across all active phases

| | |
|---|---|
| **Total SRDR Records** | 4084 |
| **Final Records with Application Domains** | 569 |
| **CSCI-Sized Records** | 447 |
| **Records With Design, Code, Test Hrs** | 377 |

**Normalization removes over 90% of the records from the dataset**

Approved for Public Release
18-MDA-9602 (23 Apr 18)

- **To use SRDR Data in Rayleigh analysis, we needed data with Rayleigh metrics populated**
  - Records need to have Peak Staff
  - Records need to have Development Months (Duration)
- **Duration calculated in months : Maximum Date (DCTI) – Minimum Date (DCTI)**

| Normalized Dataset Records | 377 |
|---|---|
| Has Peak Staff | 374 |
| Has Development Months | 373 |

- **For the remaining dataset records, solve for $t_d$**

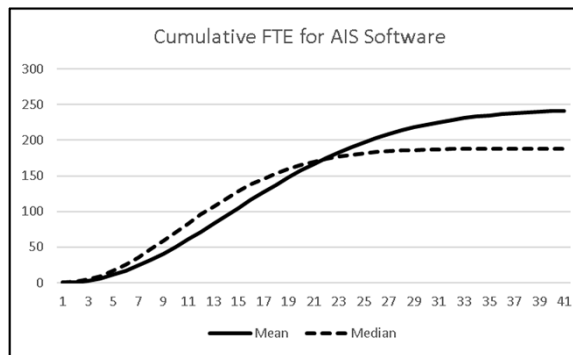$$t_d = \frac{K\,(0.6065)}{FTE_{MAX}}$$
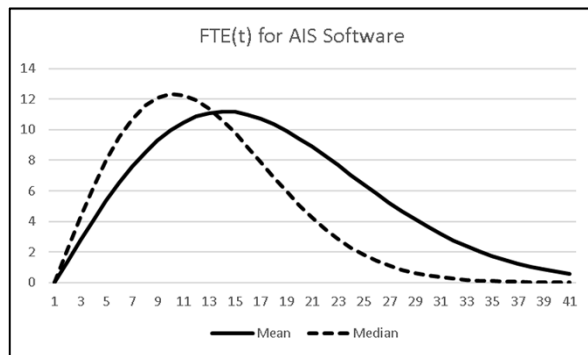
**$FTE_{MAX}$ = Peak Staff**



|  | K | Duration | td |
|---|---|---|---|
| Mean | 311.76 | 38.8 | 12.73 |
| Median | 175.90 | 32 | 9.36 |

# SRDR Data Rayleigh Curves by Super-Domain

Automated Information Systems (AIS) Software – Mission Planning, Custom Automated Information Systems, Enterprise Information Systems, and Enterprise Service Systems



| 32 records | K | Duration | td |
|---|---|---|---|
| Mean | 243.89 | 34.7 | 13.20 |
| Median | 188.25 | 24.5 | 9.25 |

Engineering Software – System Software, Process Control, Scientific/Simulation, Test/Measurement/Diagnostic Equipment
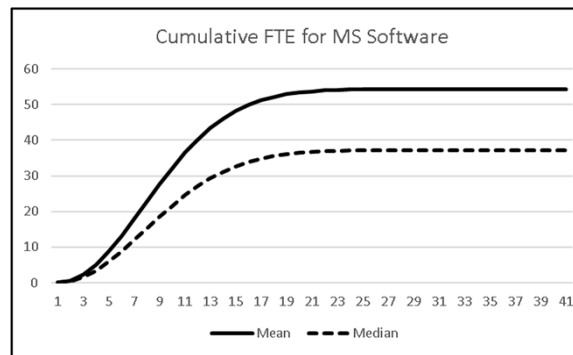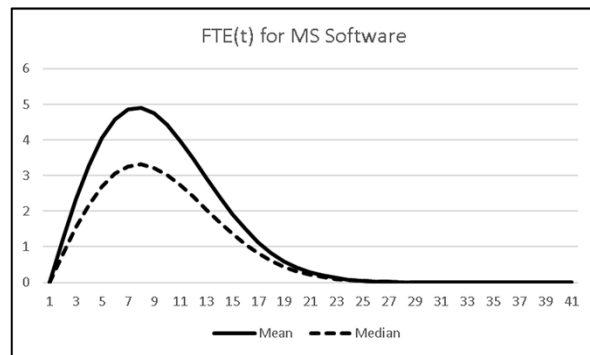


| 72 records | K | Duration | td |
|---|---|---|---|
| Mean | 270.67 | 41.7 | 13.67 |
| Median | 189.20 | 30.5 | 10.79 |

Support Software – Training and Software Tools



FTE(t) for MS Software



Cumulative FTE for MS Software

| 12 records | K | Duration | td |
|---|---|---|---|
| Mean | 54.39 | 29.2 | 6.71 |
| Median | 37.27 | 27 | 6.82 |

Real-Time Software – Microcode/Firmware, Signal Processing, Vehicle Payload/Control, Command & Control, Communications



FTE(t) for RT Software



Cumulative FTE for RT Software

| 257 records | K | Duration | td |
|---|---|---|---|
| Mean | 343.75 | 38.9 | 12.68 |
| Median | 180.64 | 35 | 9.32 |

# SRDR Quality Issues – Impossible Schedules

- **Proposed Benchmark – $t_d$ percentage = $t_d$ / Total Development Time**
- **Calculated $t_d$ percentage for the 373 normalized records**
- **Some of the $t_d$ percentage values were \*above\* 100% - peak staffing date occurs after delivery!**
- **Some of the $t_d$ percentage values were above 75% - peak staffing occurs close to delivery**
- **Calculated Maximum Load for each record**
  - Maximum Load = Peak Staff * Total Development Time
  - Records with more Development Hours than Maximum Load are in an Impossible Region

- **Removal of 38 Impossible Region records**
- **New dataset is 335 records**



Impossible Region Records
- 10% Impossible
- 90% Possible

Td Percentages
- 3% Over 100%
- 3% 75-100%
- 11% 50-70%
- 83% Under 50%

The SRDR Working Group addressed Impossible Schedule data and has added an indicator field in the SRDR database
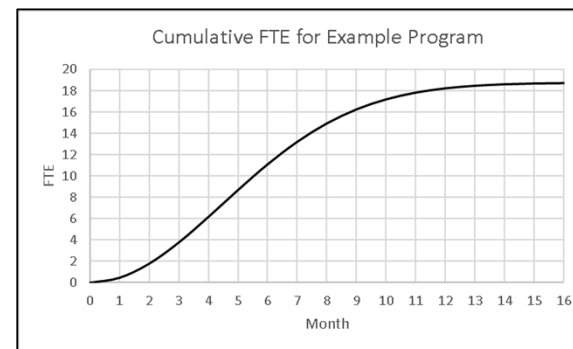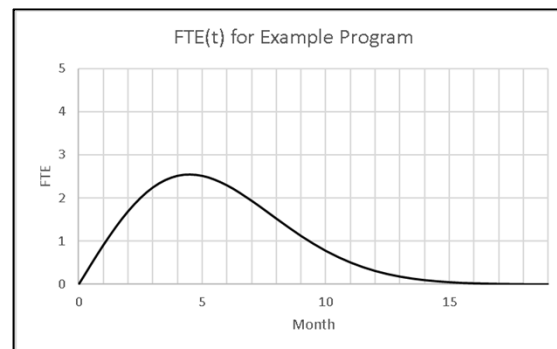
# Rayleigh Peak Staffing Percentage Benchmarks

- **Mean and median $t_d$ percentage values for the remaining 335 records calculated**

| | Td Percentage | | |
|---|---|---|---|
| | Records | Mean | Median |
| **All** | 335 | 28% | 28% |
| **AIS** | 28 | 34% | 35% |
| **Engineering** | 65 | 31% | 30% |
| **Support** | 11 | 24% | 18% |
| **Real-Time** | 231 | 27% | 26% |

> **Outside of Support Software, median and mean $t_d$ percentage values are very similar – either can be the "average"**

- **Example: Generic software program with 3000 hours, planned development schedule of 16 months**

  - **$K = 3000 / 160$ (average hrs per man-month) $= 18.75$ man months**

  - **$t_d = 16$ months $* 28\% = 4.48$ months**

  - **$FTE_{MAX} = (K * 0.6065) / t_d = 2.54$ FTE**



FTE(t) for Example Program



Cumulative FTE for Example Program

# Rayleigh as Weibull

- **Rayleigh is actually a special case of the Weibull distribution**

$$f(x; \lambda, k) = \frac{k}{\lambda}\left(\frac{x}{\lambda}\right)^{k-1} e^{-\left(\frac{x}{\lambda}\right)^k}$$

$$f(x; \lambda, k) = 1 - e^{-\left(\frac{x}{\lambda}\right)^k}$$

where k = 2 and the scale parameter ($\lambda$) = $\sqrt{2}\sigma$

- **Weibull is native to MS-Excel – easier to use than Rayleigh**
- **"=WEIBULL.DIST(x, alpha, beta, cumulative)"**
  - x = time t
  - alpha = 2
  - beta = SQRT(2) * $t_d$
  - cumulative = {FALSE (for PDF), TRUE (for CDF)}
- **Example: A generic software program estimated at 5000 hours with a duration of 16 months**
  - Use 28% for $t_d$ percentage – $t_d = \sqrt{2}*(0.28)$TDEV = 6.336
  - PDF at time x: "=5000 * WEIBULL.DIST (x, 2, 6.336, FALSE)"

Approved for Public Release
18-MDA-9602 (23 Apr 18)

# Boehm and Stutzke Rayleigh Normalizations

- **Dr. Barry Boehm suggested that pure Rayleigh is inaccurate as no project starts with zero staff and Rayleigh starts at the origin**
- **Boehm suggested using only the portion of Rayleigh from 0.3$t_d$ to 1.7$t_d$**

$$FTE(t) = K * \left( \frac{0.15 * TDEV + 0.7 * t}{0.25 * TDEV^2} \right) e^{\left( \frac{-(0.15 * TDEV + 0.7 * t)^2}{0.5 * TDEV^2} \right)}$$

- **Dick Stutzke identified that this equation was not fully normalized for the new endpoints**
- **Stutzke normalized Boehm's formula by dividing results by 1.029**

**NOTE: Because of the defined endpoints of 0.3$t_d$ and 1.7$t_d$, $t_d$ will always be 50% of the development duration using this normalization**

**Do \*NOT\* use the $t_d$ benchmarks here**

# Strickland Weibull Rayleigh Normalizations

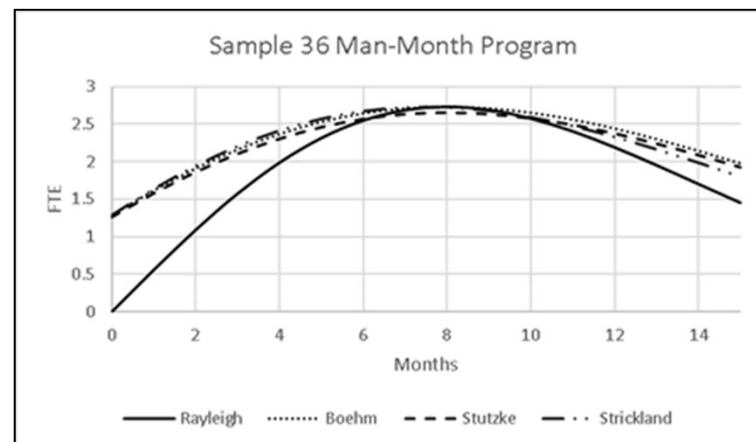- **Difference between the endpoints ($0.3t_d$ and $1.7t_d$) need to be spread over the development duration evenly**
- **Generic transformation for t:**

$$Mult(t) = 0.3 + \frac{1.4t}{(TDEV - 1)}$$

- **Expressed in Excel-like Weibull expression:**

WEIBULL.DIST (Mult(t)*$t_d$, 2, SQRT(2)*$t_d$, FALSE)

- **Example: 36 man-month program, 16 month duration**



Sample 36 Man-Month Program

**CP – Cost Analytics and Parametric Estimation Directorate**

- **Analysts don't phase using continuous distributions, would rather have hours/staffing by month**
- **Rayleigh doesn't have to start in t=0, would start in t=1**
- **Boehm / Stutzke / Strickland normalizations start in t=0 but expressed as first month**
- **Example program as a discrete distribution:**



Sample 36 Man-Month Program

**Sum of Discrete Distributions**

| Rayleigh | Boehm | Stutzke | Strickland |
| --- | --- | --- | --- |
| 31.7 | 36.7 | 35.7 | 36.2 |

**The Weibull distribution comes closest to the total in discrete calculations, but only when $t_d$ = 50% duration**

# Staffing Profiles Using Rayleigh

- **Without specific Earned Value data or effort/staffing linked in the SRDRs, we can't determine if the benchmark $t_d$ percentage or Boehm normalization is a better representation**

- **Boehm – it makes sense that staffing peaks in the middle of SW development (Code and Unit Testing)**

- **Rayleigh – the shift in paradigm in SW development is for more effort and time to be spent in design and architecture**

- **Which should you use? Either, both are better than Uniform or Triangular**

- **MDA has developed an Excel worksheet that calculates Boehm, Stutzke, and Strickland monthly phasing given inputs – used by DAU in SW Estimation course**

| Hrs | MM | Duration | Months Y1 | Total Cost ($M) |
|---|---|---|---|---|
| 5000 | 32.9 | 16 | 12 | |

| MM Distr. | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Boehm | 1.2 | 1.5 | 1.7 | 2.0 | 2.2 | 2.3 | 2.4 | 2.5 | 2.5 | 2.5 |
| | Stutzke | 1.1 | 1.4 | 1.7 | 1.9 | 2.1 | 2.2 | 2.3 | 2.4 | 2.4 | 2.4 |
| | Strickland | 1.2 | 1.5 | 1.8 | 2.0 | 2.2 | 2.3 | 2.4 | 2.5 | 2.5 | 2.4 |

| MM Per Yr | | Yr1 | Yr2 | Yr3 | Yr4 | Yr5 |
|---|---|---|---|---|---|---|
| | Boehm | 25.5 | 8.1 | 0.0 | 0.0 | 0.0 |
| | Stutzke | 24.8 | 7.9 | 0.0 | 0.0 | 0.0 |
| | Strickland | 25.5 | 7.6 | 0.0 | 0.0 | 0.0 |

# New SRDR Functionality

- **In 2017, a new SRDR DID was developed by the SRDR Working Group to address necessary changes in data fidelity and stringency**

- **The new SW Development DID includes tasking for the contractor to report SW Development hours by activity, by month, by CSCI, by build**

- **Collections of data at this level can produce staffing profile curves for software and test Rayleigh Curve metrics**
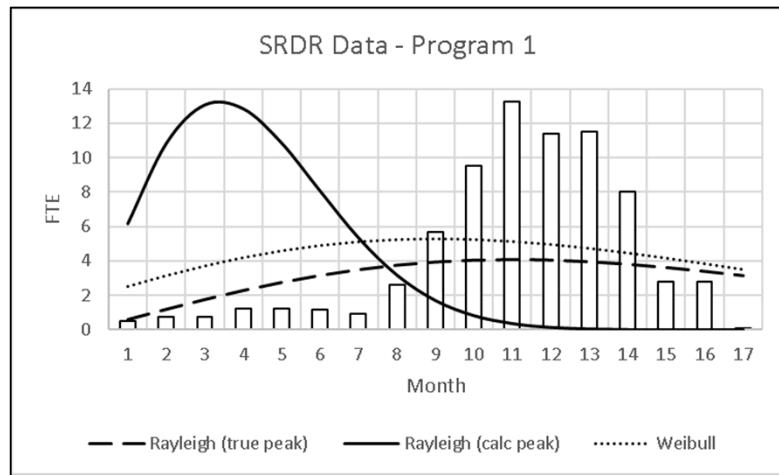
| Prime Contractor SECTION 3.4.1.1 Hours | | | | M0 YYYYMMDD | M1 YYYYMMDD | M2 YYYYMMDD | M3 YYYYMMDD | M4 YYYYMMDD | M5 YYYYMMDD |
|---|---|---|---|---|---|---|---|---|---|
| WBS Element Code | WBS Element Name | Activity ID | Activity Name | | | | | | |
| 1.1.2.2.2 | Software Release 1 | | | | | | | | |
| 1.1.2.2.2.1 | Software Release 1 CSCI 1 | | | | | | | | |
| 1.1.2.2.2.2 | Software Release 1 CSCI 1 | X | Contractor-Defined Activity X | | | | | | |
| 1.1.2.2.2.3 | Software Release 1 CSCI 1 | Y | Contractor-Defined Activity Y | | | | | | |
| 1.1.2.2.2.4 | Software Release 1 CSCI 1 | Z | Contractor-Defined Activity Z | | | | | | |
| 1.1.2.2.2.2 | Software Release 1 CSCI 2 | | | | | | | | |
| 1.1.2.2.2.2 | Software Release 1 CSCI 2 | X | Contractor-Defined Activity X | | | | | | |
| 1.1.2.2.2.2 | Software Release 1 CSCI 2 | Y | Contractor-Defined Activity Y | | | | | | |
| 1.1.2.2.2.2 | Software Release 1 CSCI 2 | Z | Contractor-Defined Activity Z | | | | | | |
| 1.1.2.2.2.3 | Software Release 1 CSCI n | | | | | | | | |
| 1.1.2.2.2.3 | Software Release 1 CSCI n | X | Contractor-Defined Activity X | | | | | | |



FTE(t) for Example Program

# SRDR Results

SRDR Data - Program 1

| Duration | 17 |
|---|---|
| K (mm) | 74.0 |
| td | 11 |
| td% | 65% |
| td (calc) | 3.4 |



SRDR Data - Program 2

| Duration | 17 |
|---|---|
| K(mm) | 264.8 |
| td | 10 |
| td% | 59% |
| td (calc) | 4.2 |

- Calculation of Rayleigh staffing curves using SRDR Final data
- Duration, K, and $t_d$ (true peak) are known, $t_d$ (calculated) is calculated given duration and K
- Program 1 – almost an inverse Rayleigh, peaked at 65% program completion
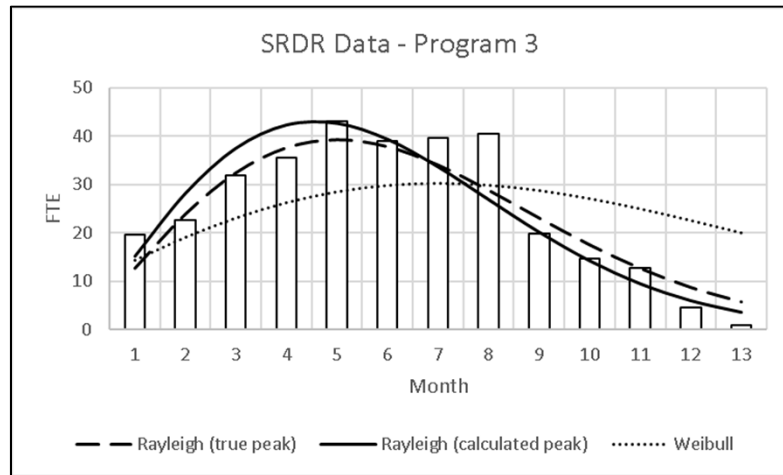- Program 2 – closer behavior to Weibull at beginning, but a steep drop-off late

# SRDR Results

SRDR Data - Program 3

| Duration | 13 |
|----------|------|
| K(mm) | 324.5 |
| td | 5 |
| td% | 38% |
| td (calc) | 4.6 |



SRDR Data - Program 4

| Duration | 5 |
|----------|------|
| K(mm) | 26.9 |
| td | 1 |
| td% | 20% |
| td (calc) | 1.3 |

- Program 3 – good match for Rayleigh, especially the calculated peak
- Program 4 – short duration, but good match for Rayleigh with a calculated peak
- Limited sample size, but Rayleigh is matching with a few programs; Weibull is not

# Future Research and Conclusions

- Work with new SRDR submissions to identify staffing profiles for software development
- Utilize Rayleigh calculators to validate SRDR submissions

- The Rayleigh curve is an acceptable time-phasing distribution for software development

- Benchmarks using real, normalized data are available

- Using the Weibull distribution as a proxy for Rayleigh yields accurate results to a Boehm-normalized curve

- Tools and benchmarks are available to help cost estimators address time-phasing in software development

# Questions